# Chapter 1

# Introduction

The past decade social media has become ever more important in our daily lives. Nowadays most of us have, beside a real world life, a virtual life on social media. This increased engagement on social media may have a huge impact on the formation of opinions, not only on the individual scale, but also on a broader, societal scale. It is known that individuals not only form their opinions through self-reflection, but also through interactions with other people and with their surroundings. The broad range of interactions people undergo on social media with people from all over the world can thus not be underestimated in the process of opinion formation. The understanding of the role of social media on the emergence of polarization and extremism in society is of uttermost importance.

One important feature of opinion dynamics on social media is the use of algorithmic personalization. Social media allow for simultaneous or asynchronous interactions between people without geographical constraints and thus allow for the spread of information at a faster pace and at an unprecedented scale [13]. People are, however, still bound by temporal and cognitive constraints. In order to assure a more pleasant/convenient experience on their platforms, social media use algorithms to order/filter the posts on an individuals time line according to what might be relevant to that individual [13]. It is thus important to know the possible effects of these filtering algorithms on the evolution of opinions in the population.

It is also known that some people change their opinions more easily than others, that is, once in a while you encounter somebody that is really stubborn and persistent towards its own opinion. One of the aims of this thesis is to design toy models of opinion dynamics with stubborn actors on theoretical and real-world social networks to analyze the interplay between individual resistance to change and the network structure in the evolution of two competing opinions.

Opinion dynamics models have two important layers. On one hand, we introduce social networks to describe the underlying structure of social interactions, while on the other hand we need appropriate opinion dynamics/formation models to reproduce the opinion dynamics and opinion formation processes found in real life.

## 1.1   Complex networks

Complex networks have become more and more popular for analyzing complex, dynamical systems. They are used in fields such as physics, economics, social sciences, biology, etc [5]. These different fields are very diverge, but they have at least one common ground: they often deal with a large number of variables/components that interact with each other. In other words, in all these fields one encounters complex systems, systems where it is not possible to predict collective behavior based on the properties of the individual components alone [7]. Complex systems often display phenomena such as non-linearity, emergence, spontaneous order,... One of the tools to deal with these complex systems and to give us more insight in the possible underlying structures are complex networks. A network, also often called a graph, is a structure composed of nodes or vertices and a set of links (edges) that indicate the interactions between the nodes [7]. Representing/modeling a complex system as a graph makes the system appear more simple and tractable, while it still includes the non-linearity of these systems. One can find the language to describe networks in mathematical graph theory. However, complex systems in real life situations often deal with a huge number of components, so the use of statistical and high-performance computing tools is inevitable [7]. One could argue that the study of complex networks lies somewhere at the intersection of mathematical graph theory and statistical physics [6].

The advantage of working with network models is that they reduce the level of complexity encountered in the real world, so that one can treat these systems in a more practical way. However we do want our models to display properties similar to the ones seen in real systems [7]. Since many real systems are not static but evolve in time, this means that we don't only need to deal with static networks but also with temporal networks. Temporal networks are networks where the edges are not always active, but instead become active for some periods of time [8]. In temporal networks time becomes an explicit element in the network representation [8]. Static networks have been widely studied and are often convenient for their analytic tractability, whereas temporal networks are, in some cases, more realistic [7]. Some other important concepts in network theory are the degree distribution, clustering/community structure, connectivity, etc. These concepts will

be explained in depth in Chapter 2, Section 2.1, but let us already anticipate a bit on the case of social networks. It is found that many real life social systems have a power law degree distribution [12]. This heterogeneous or scale free degree distribution represents one of the three general properties of social networks. The other two are short distances, also referred to as small-world phenomenon, and high clustering [12]. Ideally, our theoretical network should exhibit these three properties. Some theoretical networks, that are thoroughly studied, do not possess all of these three properties. However, they might still be used, because of their simplicity and ability to produce analytic results. When using these models one must always bear in mind their limitations to reproduce some properties encountered in real social systems.

One of the goals of this thesis is to investigate/determine the role of the underlying network structure on the formation and dynamics of opinions in the system. Some examples of theoretical network models used in this thesis are the Erdős-Rényi network, the stochastic block model, the Watts-Strogatz model, etc. These will be explained thoroughly in Chapter 2, Section 2.2.

## 1.2   Opinion dynamics

As stated before, the formation of an individual's opinion is the result of the interplay of many factors such as self-reflection, peer pressure, the individual's personality (eg. stubbornness), the information someone is exposed to, etc. Opinion dynamics models should try to capture these complex processes in a simplified way. Several models have been developed ranging from simple binary models to more complex, continuous approaches [15]. The basic idea of all opinion dynamic models is that the nodes or agents in a social network have a variable that represents their opinion and that is updated according to some predefined rules. These models are obviously a simplification of real world opinion dynamics. It is however shown that they do display a lot of aspects of real opinion formation such as agreement, transitions between order (consensus) and disorder (fragmentation), polarization, formation of echo chambers (clusters of people with the same opinion), etc [15].

The opinion of the agents in the model can be either discrete or continuous. This thesis deals with models where each agent can have one of the two opinions A or B and thus only deals with the case of discrete opinions. This might come across as a huge simplification of the real world complexity of opinion formation, but also in real life people often have to choose between two competing opinions (eg. republicans or democrats, being left-minded or right-minded, cat or dog, renting or buying a house,...).

The rules that determine how an agent updates his or her opinion can include many different aspects. They can, for example, be as simple as a majority rule (i.e. the majority model: if the majority of your neighbors have a certain opinion, you adopt that opinion as well) or can include a more probabilistic way of updating. One can also include a resistance parameter that determines the hesitation of an agent to change to a new opinion and/or a parameter that determines the influence of a node on others.

Since this thesis deals with the particular case of opinion dynamics on on-line social platforms, it is important to define/set up different filtering algorithms. The filtering algorithms that real social media companies use, are corporate secrets, but we know that there are three main principles of content curation: popularity, semantic and collaborative filtering [13]. Here popularity filtering refers to the practice of promoting content that is popular across the platform; semantic filtering means that post similar to previous consumed posts are recommended and collaborative filtering suggests posts that are similar to the ones our friends consume [13]. In Chapter 3 a detailed description of the filtering algorithms used in this thesis will be given.

## 1.3   Statistical and social physics

Another important tool in the study of opinion dynamics/formations in large groups of people is the use of statistical physics. This branch of physics offers a framework that relates microscopic properties of atoms, molecules,... to macroscopic observed behavior and has a wide field of applications inside and outside the world of physics. The observation that large number of people display collective, 'macroscopic' behavior begged for the use of the concepts and insights developed in statistical physics [15]. The application of the theory of statistical physics on social phenomena is referred to as social physics or sociophysics. The 'microscopic' constituents are now individual humans who interact with a limited number of other individuals and in that way form complex, 'macroscopic' groups such as human societies. These 'macroscopic' groups display stunning regularities, transitions from disorder to order, the emergence of consensus, universality, etc. The statistical physics approach of these social systems tries to explain the found regularities at large scale as collective effects of the interactions among these individuals [15].

## 1.4   This thesis: hypotheses and goals

As said before, this thesis will try to investigate the interplay between the social network structure and the individual resistance to change in the evolution of two competing opinions. So on one hand, we want to determine the influence of different network structures on the formation and evolution of opinions, whereas, on the other hand, we would like to investigate the effect of stubborn agents in the network. Since this thesis deals with opinion dynamics on social media, the effect of filtering algorithms must also be included.

In one of their previous papers ([13]) prof. Nicola Perra and prof. Luis E.C. Rocha investigated both the effect of different filtering algorithms and different network topologies on the formation and evolution of two competing opinions. They also considered the effect of nudging (meaning that one opinion is pushed to all agents in the network). Their main findings were that the algorithmic filtering could not break the status quo when the prevalence of opinions was equally distributed in the population and that topological correlations (such as high clustering coefficient) could result in the formation of echo chambers. In the case of nudging, they found that the population opinion moved to the nudged opinion relatively fast, even in the case of small nudging. They, however, did not investigate the effect of stubborn agents in the network, nor did they investigate networks with a community structure. This thesis will build on their model and incorporate these new features. We will compare networks that have a community structure and relatively low clustering to networks without community structure, but with a high clustering coefficient (a detailed description of concepts like clustering coefficient and community structure will be given in Chapter 2, Section 2.1). We want to investigate whether networks with community structure are able to form echo chambers, such as one can observe in networks with a high clustering coefficient and, if so, if one can observe convergence to one opinion inside the communities. On the other hand we want to determine the impact of stubborn agents on the formation and evolution of opinions in the network.

# Chapter 2

# Networks

## 2.1 Definition and network measurements

Networks or graphs are conceptually very simple and flexible objects that are made of a set of nodes/vertices $V$ and a set of edges/links $E$, where the elements of $E$ determine connections between elements of $V$ [4]

$$E \subseteq \{\{x, y\}|x, y \in V\}. \tag{2.1}$$

If self-loops are not allowed, we need the following condition on the elements of $E$

$$E \subseteq \{\{x, y\}|x, y \in V \text{ and } x \neq y\}. \tag{2.2}$$

A graph with no self-loops is a simple graph. The nodes $x$ and $y$ of an edge $\{x, y\}$ are called the endpoints of the edge. It is possible that nodes are not joined by any edge, such nodes are called isolated. If two or more edges have the same pair of endpoints, the graph is called a multi-graph. This thesis is not concerned with multi-graphs, nor are self-loops allowed. Graphs in which the edges have an orientation are called directed graphs. If the edges have weights $w_{ij}$, one speaks of a weighted graph.
A graph can be represented in different ways. One of the most common ways is by use of the adjacency matrix.

### 2.1.1    Adjacency matrix

The adjacency matrix $\mathbf{A}$ is a $N \times N$ matrix (for a graph with $N$ nodes) where the elements indicate whether two nodes are connected by an edge or not [7]

$$A_{ij} = \begin{cases} 1 & \text{if nodes } i \text{ and } j \text{ are connected} \\ 0 & \text{otherwise} \end{cases} \tag{2.3}$$

For a simple, undirected graph, the adjacency matrix is symmetric ($A_{ij} = A_{ji}$) with zeros on the diagonal ($A_{ii} = 0$). Directed graphs have asymmetric adjacency matrices. If we are dealing with weighted graphs, the elements of the adjacency matrix are the weights $w_{ij}$ of the edges

$$A_{ij} = \begin{cases} w_{ij} & \text{if nodes } i \text{ and } j \text{ are connected} \\ 0 & \text{otherwise} \end{cases} \tag{2.4}$$

with, generally, $0 \leqslant w_{ij} \leqslant 1$ [7].

### 2.1.2    Degree and degree distribution

The adjacency matrix contains a lot of information such as the degree of a node. The degree $k_i$ of a node $i$ is the number of edges attached to the node $i$ or, in others words, $k_i$ represents the number of nearest neighbors of the node $i$. The degree can be obtained from the adjacency matrix in the following way [7]

$$k_i = \sum_{j=1}^{N} A_{ij}. \tag{2.5}$$

In the case of a directed graph one can differentiate between the number of incoming edges $k_i^{\text{in}}$ and the number of outgoing edges $k_i^{\text{out}}$ of the node $i$. These are called the in-degree and out-degree respectively and are defined as [7]

$$k_i^{\text{in}} = \sum_{j=1}^{N} A_{ji} \quad \text{and} \quad k_i^{\text{out}} = \sum_{j=1}^{N} A_{ij}. \tag{2.6}$$

The total degree of node $i$ is then $k_i = k_i^{\text{in}} + k_i^{\text{out}}$. For weighted graphs the degree is easily generalized to the weighted degree $s_i$, often called strength [9]

$$s_i = \sum_{j=1}^{N} w_{ij}. \tag{2.7}$$

From the degree of the nodes in the network one can construct the degree distribution. The degree distribution $P(k)$ represents the fraction of nodes with a degree $k$ or, in other words, the degree distribution gives the probability that a randomly selected node has a degree $k$. The average degree can be obtained from the degree distribution [7]

$$\langle k \rangle = \frac{1}{N} \sum_{i=1}^{N} k_i = \sum_k k P(k) \tag{2.8}$$

where $N$ is the total number of nodes in the network. The degree distribution also allows us to classify networks. The two most important classes are homogeneous and heterogeneous networks. Homogeneous networks have a bell curved degree distribution, eg. a Poisson distribution. In this case, most of the nodes have a degree closed to the average degree $k$. Heterogeneous networks, on the other hand, have a power-law degree distribution, $P(k) \sim k^{-\gamma}$. They are also referred to as scale free networks, since they do not posses a characteristic length scale or, in this case, the average degree isn't a characteristic scale for the network [7].

### 2.1.3 Connectivity

### 2.1.4 Clustering coefficient and transitivity

The clustering coefficient is another important measure of network topology. It is a measure for the number of triangles in a network and determines the connectivity in the neighborhood of a node $i$: if a node $i$ has a high clustering coefficient, its neighbors are likely to be directly connected to each other [11]. Here a triangle is defined as a loop of length three, this is a sequence of nodes $x, y, z, x$ such that $\{x, y\}, \{y, z\}$ and $\{z, x\}$ are edges of the network. The clustering coefficient is thus a way to measure the degree to which nodes in a network tend to cluster [11]. The local clustering coefficient of a node $i$ is defined as

$$C_i = \frac{2n_i}{k_i(k_i - 1)} \tag{2.9}$$

where $n_i$ is the number of edges that actually exist between the nodes in the neighborhood of $i$ and $k_i(k_i - 1)/2$ is the maximum number of edges that could exist between them (note that this expression is only valid for an undirected graph, since for a directed graph $e_{ij} \neq e_{ji}$ and we have $k_i(k_i - 1)$ possible edges between the neighbors of node i). The average clustering coefficient is then given as the average of the local clustering coefficients

$$\bar{C} = \frac{1}{N} \sum_{i=1}^{N} C_i \tag{2.10}$$

where $N$ denotes, once again, the number of nodes in the network. The clustering coefficient is, beside a measure for the connectivity in the network, also linked to the robustness, or resilience against random damage, of the network [10] [11].

A measure that is closely related to the clustering coefficient is the transitivity $T$, defined as (valid for undirected, unweighted networks) [6]

$$T = \frac{3 \times \text{number of triangles in the network}}{\text{number of connected triples of nodes in the network}}. \tag{2.11}$$

A connected triple is defined as a set of three nodes with at least two edges between them, so that each node can be reached from the other two (either directly or indirectly). The factor three arises from the fact that each triangle contributes to three different connected triples in the network: one centered at each node in the triangle [6]. Let us denote the number of triangles as $N_\Delta$ and the number of connected triples as $N_3$. These two numbers can be obtained from the adjacency matrix in the following way [6]

$$N_\Delta = \sum_{k>j>i} a_{ij} a_{ik} a_{jk} \tag{2.12}$$

$$N_3 = \sum_{k>j>i} (a_{ij} a_{ik} + a_{ji} a_{jk} + a_{ki} a_{kj}) \tag{2.13}$$

where $a_{ij}$ are the elements of the adjacency matrix. It is possible to define the transitivity of one node as [6]

$$T_i = \frac{N_\Delta(i)}{N_3(i)} \tag{2.14}$$

where $N_\Delta(i)$ represents the number of triangles that involve node $i$ and $N_3(i)$ is the number of connected triples with $i$ as the central node [6]

$$N_\Delta(i) = \sum_{k>j} a_{ij} a_{ik} a_{jk} \tag{2.15}$$

$$N_3(i) = \sum_{k>j} a_{ij} a_{ik} \tag{2.16}$$

It is not too hard to see that $N_\Delta(i)$ counts the number of edges between the neighbors of $i$ and that $N_3(i)$ is equal to $k_i(k_i - 1)/2$, where $k_i$ is the degree of node $i$. It is thus obvious that the local clustering coefficient $C_i$ (Eq. (2.9)) of a node and the transitivity $T_i$ (Eq. (2.14)) of a node define the same quantity [6]. In the case of the average clustering coefficient $\bar{C}$ and the global transitivity $T$ this is however not the case. The difference between the two definitions is that the average of

Equation (2.11) gives the same weight to each triangle in the network, whereas Equation (2.10) gives the same weight to each node. This may lead to slightly different values since nodes with a higher degree may possibly be involved in a higher number of triangles than nodes with a lower degree [6].

### 2.1.5 Network communities

Network communities is another feature that (complex) networks may posses. A community inside a network is loosely defined as a set of nodes that are more densely connected to nodes inside that set than to the other nodes in the network [14]. More strictly, a community is defined based on two hypotheses: the connectedness hypothesis and the density hypothesis [3]. In short, the connectedness hypothesis means that each member of a community should be reached through each other member of the same community [3]. The density hypothesis implies that nodes inside a community are more likely to be linked to other nodes inside that community than to nodes outside the community [3]. The density hypothesis narrows what could be considered a community, but it doesn't uniquely define it. Several community definitions are consistent with the density hypothesis. Let us consider three possible definitions: maximum cliques, strong community and weak community [3].

**Maximum cliques**

A clique is a complete subgraph, where a complete subgraph is defined as a set of nodes of the network where each node in the set is directly connected to all the others nodes in the same set [3]. A community based on this definition would then be the largest clique in the network. This definition of community might, however, be too restrictive and, beside that, large cliques don't appear very frequent in networks [3].

**Strong communities**

A strong community is defined such that each node inside the community has more links to other nodes inside the same community than to nodes outside the community [3]. Let us denote a community as $C$ and let $i$ be a node inside the community $C$. If we define the internal degree $k_i^{int}(C)$ as the number of links of node $i$ with other nodes in $C$ and the external degree $k_i^{ext}(C)$ as the number of links of node $i$ with nodes outside $C$, then we have the following condition for $C$ to be a strong community [3]

$$k_i^{int}(C) > k_i^{ext}(C) \qquad \forall i \in C. \tag{2.17}$$

**Weak communities**

A weak community is a community where the sum of the internal degree of all the nodes in the community exceeds the sum of the external degree of all the nodes in the community [3]. We thus have the following condition for a weak community $C$ [3]

$$\sum_{i \in C} k_i^{int}(C) > \sum_{i \in C} k_i^{ext}(C).$$

(2.18)

## 2.2   Network models

In order to study the topological structures of real-world networks, several theoretical network models have been developed. These theoretical models usually have a simpler representation and well known properties that can be derived mathematically. They are widely studied and used extensively. The study of complex networks relies heavily on the knowledge and understanding of these models. In this section we will review some of the most important ones.

### 2.2.1   The Erdős-Rényi network

In 1959 the two mathematicians Paul Erdős and Alfréd Rényi proposed a model to generate random networks [6]. Independently of them, Solomonoff and Rapoport already proposed the model in 1951. The model constructs random networks with $N$ nodes in the following way [1]

1. Start with $N$ disconnected nodes,

2. For each pair of nodes, connect them with a predefined probability $p$.

If $p = 1$, we obtain a network with the maximum number of edges $N(N-1)/2$. This is called the fully connected network.
The probability that a node is connected to $k$ other nodes and not to $N - k$ others is given by

$$p^k (1-p)^{N-1-k}.$$

(2.19)

The number of ways to choose the $k$ nodes among the $N - 1$ possible candidates is given by the binomial coefficient

$$\binom{N-1}{k} = \frac{(N-1)!}{k!(N-1-k)!}$$

(2.20)

and the probability of being connected to exactly $k$ other nodes becomes

$$p_k = \binom{N-1}{k} p^k (1-p)^{N-1-k}. \tag{2.21}$$

We thus derived the degree distribution of a random network. From this we can calculate the average degree of the network. The expected mean degree is given by [16]

$$
\begin{aligned}
\langle k \rangle &= \sum_{k=0}^{N-1} k \, p_k \\
&= \sum_{k=0}^{N-1} k \binom{N-1}{k} p^k (1-p)^{N-1-k}.
\end{aligned}
\tag{2.22}
$$

In order to simplify this equation we can make use of the following formula [16]

$$(p+q)^n = \sum_{k=0}^{n} \binom{n}{k} p^k q^{n-k}. \tag{2.23}$$

Differentiating both sides of this equation gives us [16]

$$
\begin{aligned}
n(p+q)^{n-1} &= \sum_{k=0}^{n} \binom{n}{k} k p^{k-1} q^{n-k} \\
&= \frac{1}{p} \sum_{k=0}^{n} \binom{n}{k} k p^k q^{n-k}.
\end{aligned}
\tag{2.24}
$$

Finally, by substituting $q = 1 - p$, we get [16]

$$
\begin{aligned}
n(p+(1-p))^{n-1} &= \frac{1}{p} \sum_{k=0}^{n} \binom{n}{k} k p^k (1-p)^{n-k} \\
np &= \sum_{k=0}^{n} \binom{n}{k} k p^k (1-p)^{n-k}
\end{aligned}
\tag{2.25}
$$

where the right hand side is equal to Equation (2.22) with $n = N - 1$ and we thus obtain $\langle k \rangle = (N-1)p$. [2] large N

# Chapter 3

# Opinion dynamics